

CS498 AMO

Name: Xin Qu, ID: xinq2

Part 1 Accuracies

Setup	Cross_validation Accuracy
Unprocessed data	75.45%
0-value elements ignored	75.83%

1&2. Calculation of distribution parameters and Naive Bayes predictions

```
In [4]: 1 def norm_prob(x, mean_y, var_y): ### var = stdv**2
2         ###return P(x|y)
3         tmp = np.exp(-(x - mean_y) ** 2 / (2 * var_y))
4         return 1 / ((2 * np.pi * var_y) ** 0.5) * tmp

In [5]: 1 def class_prob(xi, data):
2         ###return class probability
3         prob = {0: 1, 1: 1}
4         for i in range(len(xi) - 1): ### xi needs to drop column('class')
5             ###P(xi|y) = P(x0|y) * P(x1|y)*...
6             prob[0] = prob[0] * norm_prob(xi[i], data['mean_y'].iloc[0][i], data['var_y'].iloc[0][i])
7             prob[1] = prob[1] * norm_prob(xi[i], data['mean_y'].iloc[1][i], data['var_y'].iloc[1][i])
8         prob[0] *= data['pt']
9         prob[1] *= data['pn']
10        return prob
11 def classify(xi, data):
12     prob = class_prob(xi, data)
13     return 0 if prob[0] > prob[1] else 1

In [6]: 1 def NBclassifier(dataset, ratio):
2         accu_list = []
3         for i in range(10):
4             train, test = train_test_split(dataset, ratio)
5             class_1 = train['class'][train['class'] == 1].count()
6             class_0 = train['class'][train['class'] == 0].count()
7             class_total = train.shape[0]
8             pn = class_1 / float(class_total) ##1:negative
9             pt = class_0 / float(class_total) ## 0: positive
10            mean_y = train.groupby('class').mean()
11            var_y = train.groupby('class').var()
12            ##std_y = train.groupby('class').std()
13            cor = 0
14            for i, j in test.iterrows():
15                result = classify(j, {'mean_y': mean_y, 'var_y': var_y, 'pn':pn, 'pt':pt})
16                if bool(result) == bool(j['class']):
17                    cor += 1
18            accu_list.append(cor / float(test.shape[0]))
19        accuracy = np.sum(accu_list) / 10.0
20        return accuracy
```

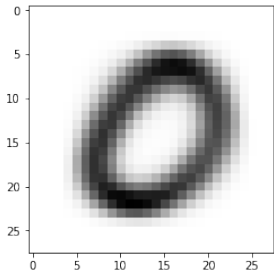
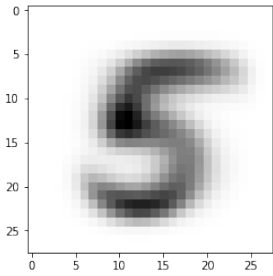
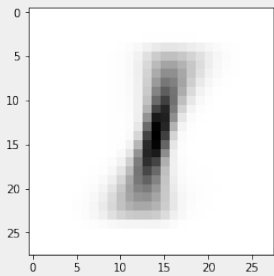
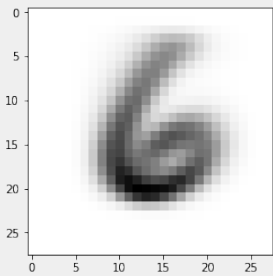
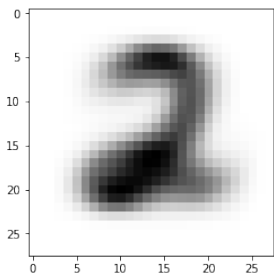
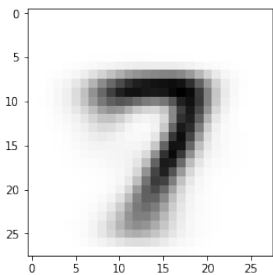
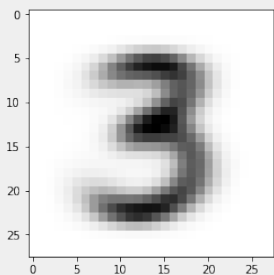
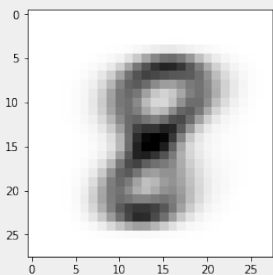
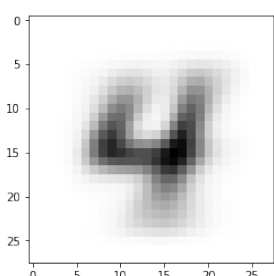
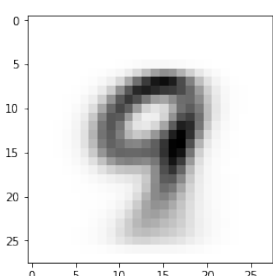
3. Test-train split code

```
In [3]: 1 import random
2         def train_test_split(dataset, ratio):
3             data = dataset.values
4             random.shuffle(data)
5             size = len(data)
6             train_size = int(size * (1 - ratio))
7             train_data = data[:train_size, :]
8             test_data = data[train_size:, :]
9             train = pd.DataFrame(train_data, index = None,
10                                 columns = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class'])
11             test = pd.DataFrame(test_data, index = None,
12                                 columns = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class'])
13             train[['preg', 'plas', 'pres', 'skin', 'test', 'age', 'class']] = train[['preg',
14                                                                                     'plas', 'pres', 'skin', 'test', 'age', 'class']].astype(int)
15             test[['preg', 'plas', 'pres', 'skin', 'test', 'age', 'class']] = test[['preg',
16                                                                                     'plas', 'pres', 'skin', 'test', 'age', 'class']].astype(int)
17             return train, test
```

Part 2 MNIST Accuracies

x	Method	Training Set Accuracy	Test Set Accuracy
1	Gaussian + untouched	77.75%	79.58%
2	Gaussian + stretched	79.75%	82.06%
3	Bernoulli + untouched	83.77%	85.17%
4	Bernoulli + stretched	82.50%	84.15%
5	10 trees + 4 depth + untouched	74.33%	76.15%
6	10 trees + 4 depth + stretched	71.62%	72.98%
7	10 trees + 16 depth + untouched	99.46%	94.89%
8	10 trees + 16 depth + stretched	99.48%	95.53%
9	30 trees + 4 depth + untouched	79.83%	80.61%
10	30 trees + 4 depth + stretched	74.20%	75.76%
11	30 trees + 16 depth + untouched	99.76%	96.20%
12	30 trees + 16 depth + stretched	99.72%	96.91%

Part 2A Digit Images

Digit		Mean Image	
0			
5			
1			
6			
2			
7			
3			
8			
4			
9			

1& 2. Calculation of Normal distribution and Bernoulli distribution parameters & 3. Calculation of Naive Bayes predictions

```

1 def Gaussiannb_predict(trainset, testset):
2     def norm_prob_cal(x, mean_y, var_y): ### var = stdv**2
3         ###return P(x|y)
4         tmp = np.exp(-(x - mean_y) ** 2 / (2 * var_y))
5         return 1 / ((2 * np.pi * var_y) ** 0.5) * tmp
6     #model_nb = GaussianNB()#model_nb.fit(trainset, train_y)
7     #mean = model_nb.theta_ ### mean values (10, 784)
8     #var = model_nb.sigma_ ### var values (10, 784)
9     #pro = model_nb.class_prior_ ###prob of each class (10)
10    from collections import Counter
11    count_dic = Counter(testset)
12    prob = []
13    for i in count_dic.keys():
14        p = count_dic[i] / float(sum(count_dic.values()))
15        prob.append(p)
16    prob = np.array(prob)
17    columns = [str(i) for i in range(trainset.shape[1])]
18    df = pd.DataFrame(trainset, index = None, columns = columns)
19    df['class'] = testset
20    mean_ = df.groupby('class').mean()
21    var_ = df.groupby('class').var()
22    pro = prob
23    mean = mean_
24    var = var_
25    pred_y = []
26    ###return predict of trainset
27    for i in range(len(trainset)):
28        norm_prob = norm_prob_cal(trainset[i], mean, var)
29        norm_prob = norm_prob.replace(np.NaN, 1)
30        norm_log_prob = np.log(norm_prob)
31        p = np.sum(norm_log_prob, axis = 1) + np.log(pro)
32        pred_y.append(np.argmax(p))
33    return pred_y

1 def Bernoulli_pred(trainset_x, trainset_y):
2     #model = BernoulliNB()
3     #model.fit(trainset, train_y)
4     #log_p_class = model.class_log_prior_
5     #log_p_feature = model.feature_log_prob_
6     columns = [str(i) for i in range(trainset_x.shape[1])]
7     df = pd.DataFrame(trainset_x, index = None, columns = columns)
8     df['class'] = trainset_y
9     count_dic = Counter(trainset_y)
10    prob = []
11    for i in count_dic.keys():
12        p = count_dic[i] / float(sum(count_dic.values()))
13        prob.append(p)
14    prob = np.array(prob)
15    log_p_class = np.log(prob)
16    class_dic = Counter(trainset_y)
17    class_count = np.array(count_dic.values())
18    feature_count = np.array(df.groupby('class').sum())
19    log_p_feature = np.log(feature_count + 1.0) - np.log((class_count + 2.0).reshape(-1, 1))
20    lpg_p_feature = np.array(log_p_feature)
21    pred = []
22    for i in range(len(trainset_x)):
23        tmp = trainset_x[i] * log_p_feature + (1 - trainset_x[i]) * np.log((1 - np.exp(log_p_feature)))
24        total_tmp = np.sum(tmp, axis = 1)
25        p = total_tmp + log_p_class
26        pred.append(np.argmax(p))
27    return pred

```

4. Training of a decision tree & 5. Calculation of a decision tree predictions

```

1 resultB = defaultdict(list)
2 for (train, test, name) in dataset:
3     for [tree, depth] in combinations:
4         model = RandomForestClassifier(n_estimators = tree, max_depth = depth)
5         model.fit(train, train_y)
6         pred_y = model.predict(test)
7         train_score = accuracy_score(map(int, train_y), model.predict(train))
8         test_score = accuracy_score(map(int, test_y), pred_y)
9         resultB[name].append((train_score, test_score))

```

Code for part 1

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 df = pd.read_csv('pima-indians-diabetes.csv', header = None, names =
4               ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class'])
5 df.head()
```

```
Out[1]:
```

	preg	plas	pres	skin	test	mass	pedi	age	class
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

```
In [2]: 1 df.shape
```

```
Out[2]: (768, 9)
```

```
In [3]: 1 import random
2 def train_test_split(dataset, ratio):
3     data = dataset.values
4     random.shuffle(data)
5     size = len(data)
6     train_size = int(size * (1 - ratio))
7     train_data = data[:train_size, :]
8     test_data = data[train_size:, :]
9     train = pd.DataFrame(train_data, index = None,
10                       columns = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class'])
11     test = pd.DataFrame(test_data, index = None,
12                       columns = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class'])
13     train[['preg', 'plas', 'pres', 'skin', 'test', 'age', 'class']] = train[['preg',
14     'plas', 'pres', 'skin', 'test', 'age', 'class']].astype(int)
15     test[['preg', 'plas', 'pres', 'skin', 'test', 'age', 'class']] = test[['preg',
16     'plas', 'pres', 'skin', 'test', 'age', 'class']].astype(int)
17     return train, test
```

```
In [4]: 1 def norm_prob(x, mean_y, var_y): ### var = stdv**2
2     ###return P(x|y)
3     tmp = np.exp(-(x - mean_y) ** 2 / (2 * var_y))
4     return 1 / ((2 * np.pi * var_y) ** 0.5) * tmp
```

```
In [5]: 1 def class_prob(xi, data):
2     ###return class probability
3     prob = {0: 1, 1: 1}
4     for i in range(len(xi) - 1): ### xi needs to drop column('class')
5         ###P(xi|y) = P(x0|y) * P(x1|y)*...
6         prob[0] = prob[0] * norm_prob(xi[i], data['mean_y'].iloc[0][i], data['var_y'].iloc[0][i])
7         prob[1] = prob[1] * norm_prob(xi[i], data['mean_y'].iloc[1][i], data['var_y'].iloc[1][i])
8         prob[0] *= data['pt']
9         prob[1] *= data['pn']
10    return prob
11 def classify(xi, data):
12     prob = class_prob(xi, data)
13     return 0 if prob[0] > prob[1] else 1
```

```
In [6]: 1 def NBclassifier(dataset, ratio):
2     accu_list = []
3     for i in range(10):
4         train, test = train_test_split(dataset, ratio)
5         class_1 = train['class'][train['class'] == 1].count()
6         class_0 = train['class'][train['class'] == 0].count()
7         class_total = train.shape[0]
8         pn = class_1 / float(class_total) ##1:negative
9         pt = class_0 / float(class_total) ## 0: positive
10        mean_y = train.groupby('class').mean()
11        var_y = train.groupby('class').var()
12        ##std_y = train.groupby('class').std()
13        cor = 0
14        for i, j in test.iterrows():
15            result = classify(j, {'mean_y': mean_y, 'var_y': var_y, 'pn':pn, 'pt':pt})
16            if bool(result) == bool(j['class']):
17                cor += 1
18        accu_list.append(cor / float(test.shape[0]))
19    accuracy = np.sum(accu_list) / 10.0
20    return accuracy
```

```
In [7]: 1 accuracy = NBclassifier(df, 0.2)
2 print('The estimate of accuracy of the classifier is %.2f' %(accuracy * 100) + '%')

The estimate of accuracy of the classifier is 75.45%
```

part 1B

```
In [8]: 1 df_B = pd.read_csv('pima-indians-diabetes.csv', header = None,
2 names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class'])
3 df_B[['pres', 'skin', 'mass', 'age']] = df_B[['pres', 'skin', 'mass', 'age']].replace(0, np.NaN)
4 df_B = df_B.dropna()
5 df_B.head()
6 df_B.shape
```

Out[8]: (537, 9)

```
In [9]: 1 accuracy_B = NBclassifier(df_B, 0.2)
```

```
In [10]: 1 print('The estimate of accuracy of the modified classifier by is %.2f' %(accuracy_B * 100) + '%')

The estimate of accuracy of the modified classifier by is 75.83%
```

Code for part 2

```
1 from mnist import MNIST
2 mndata = MNIST('/Users/xinqu/Sandbox/CS498 Applied Machine Learning/HW1/python-mnist/data')
3 train_x, train_y = mndata.load_training()
4 test_x, test_y = mndata.load_testing()
```

```
1 from sklearn.naive_bayes import GaussianNB
2 from sklearn.naive_bayes import BernoulliNB
3 from sklearn.metrics import accuracy_score
4 import matplotlib.pyplot as plt
```

```
1 from skimage.transform import resize
2 def crop_image(x):
3     hor_max, ver_max = np.max(np.where(x != 0), 1)
4     hor_min, ver_min = np.min(np.where(x != 0), 1)
5     bound = x[hor_min: hor_max, ver_min: ver_max]
6     resize_x = resize(bound, (20, 20), mode = 'constant')
7     return resize_x
```

```
1 def threshold(x):
2     y = x.copy()
3     grey_y = y[y > 0] ## remove 0 pixel
4     mid_grey = sum(grey_y) / len(grey_y)
5     y[y < mid_grey] = 0
6     y[y > 0] = 1
7     return y
```

```
1 import numpy as np
2 train_x = np.array(train_x).astype(np.uint8)
3 train_y = np.array(train_y).astype(np.uint8)
4 test_x = np.array(test_x).astype(np.uint8)
5 test_y = np.array(test_y).astype(np.uint8)
```

```
1 ##untouched image
2 ###threshold pixel
3 train_x_thre = []
4 for x in train_x:
5     train_x_thre.append(x)
6 train_x_thre = np.array([threshold(x) for x in train_x_thre])
7 test_x_thre = []
8 for x in test_x:
9     test_x_thre.append(x)
10 test_x_thre = np.array([threshold(x) for x in test_x_thre])
```

```
1 ##stretched image
2 ###crop image
3 train_x_crop = []
4 test_x_crop = []
5 for x in train_x.reshape(-1, 28, 28):
6     train_x_crop.append(np.ravel(crop_image(x)))
7 for x in test_x.reshape(-1, 28, 28):
8     test_x_crop.append(np.ravel(crop_image(x)))
9 test_x_crop = np.array(test_x_crop)
10 ###
11 train_x_crop_thre = np.array([threshold(x) for x in train_x_crop])
12 test_x_crop_thre = np.array([threshold(x) for x in test_x_crop])
```

```
1 def accuracyscore(x, y):
2     return sum(x == y) / float(len(y))
```

```

1 def Gaussiannb_predict(trainset, testset):
2     def norm_prob_cal(x, mean_y, var_y): ## var = stdv**2
3         ##return P(x|y)
4         tmp = np.exp(-(x - mean_y) ** 2 / (2 * var_y))
5         return 1 / ((2 * np.pi * var_y) ** 0.5) * tmp
6     #model_nb = GaussianNB()#model_nb.fit(trainset, train_y)
7     #mean = model_nb.theta_ ### mean values (10, 784)
8     #var = model_nb.sigma_ ### var values (10, 784)
9     #pro = model_nb.class_prior_ ###prob of each class (10)
10    from collections import Counter
11    count_dic = Counter(testset)
12    prob = []
13    for i in count_dic.keys():
14        p = count_dic[i] / float(sum(count_dic.values()))
15        prob.append(p)
16    prob = np.array(prob)
17    columns = [str(i) for i in range(trainset.shape[1])]
18    df = pd.DataFrame(trainset, index = None, columns = columns)
19    df['class'] = testset
20    mean_ = df.groupby('class').mean()
21    var_ = df.groupby('class').var()
22    pro = prob
23    mean = mean_
24    var = var_
25    pred_y = []
26    ##return predict of trainset
27    for i in range(len(trainset)):
28        norm_prob = norm_prob_cal(trainset[i], mean, var)
29        norm_prob = norm_prob.replace(np.NaN, 1)
30        norm_log_prob = np.log(norm_prob)
31        p = np.sum(norm_log_prob, axis = 1) + np.log(pro)
32        pred_y.append(np.argmax(p))
33    return pred_y

```

```

1 accuracyscore(Gaussiannb_predict(train_x_thre, train_y), train_y) ###Gaussian train set untouched

```

/anaconda2/lib/python2.7/site-packages/ipykernel_launcher.py:30: RuntimeWarning: divide by zero encountered in log

0.7775333333333333

```

1 accuracyscore(Gaussiannb_predict(train_x_crop_thre, train_y), train_y) ###Gaussian train set stretched

```

/anaconda2/lib/python2.7/site-packages/ipykernel_launcher.py:30: RuntimeWarning: divide by zero encountered in log

0.79585

```

1 accuracyscore(Gaussiannb_predict(test_x_thre, test_y), test_y) ###Gaussian test set untouched

```

0.7975

```

1 accuracyscore(Gaussiannb_predict(test_x_crop_thre, test_y), test_y) ###Gaussian test set stretched

```

0.8206

```

1 def Bernoulli_pred(trainset_x, trainset_y):
2     #model = BernoulliNB()
3     #model.fit(trainset, train_y)
4     #log_p_class = model.class_log_prior_
5     #log_p_feature = model.feature_log_prob_
6     columns = [str(i) for i in range(trainset_x.shape[1])]
7     df = pd.DataFrame(trainset_x, index = None, columns = columns)
8     df['class'] = trainset_y
9     count_dic = Counter(trainset_y)
10    prob = []
11    for i in count_dic.keys():
12        p = count_dic[i] / float(sum(count_dic.values()))
13        prob.append(p)
14    prob = np.array(prob)
15    log_p_class = np.log(prob)
16    class_dic = Counter(trainset_y)
17    class_count = np.array(count_dic.values())
18    feature_count = np.array(df.groupby('class').sum())
19    log_p_feature = np.log(feature_count + 1.0) - np.log((class_count + 2.0).reshape(-1, 1))
20    lpg_p_feature = np.array(log_p_feature)
21    pred = []
22    for i in range(len(trainset_x)):
23        tmp = trainset_x[i] * log_p_feature + (1 - trainset_x[i]) * np.log((1 - np.exp(log_p_feature)))
24        total_tmp = np.sum(tmp, axis = 1)
25        p = total_tmp + log_p_class
26        pred.append(np.argmax(p))
27    return pred

```

```

1 accuracyscore(Bernoulli_pred(train_x_thre, train_y), train_y)## Bernoulli untouched train set

```

0.8376833333333333

```

1 accuracyscore(Bernoulli_pred(test_x_thre, test_y), test_y)## Bernoulli untouched test set

```

0.8517


```
1 accuracyscore(Bernoulli_pred(train_x_crop_thre, train_y), train_y)## Bernoulli stretched train set
```

```
0.8250333333333333
```

```
1 accuracyscore(Bernoulli_pred(test_x_crop_thre, test_y), test_y)## Bernoulli stretched test set
```

```
0.8415
```

plot mean pixel values for Normal Distribution

```
1 columns = [str(i) for i in range(train_x_thre.shape[1])]
2 df = pd.DataFrame(train_x_thre, index = None, columns = columns)
3 df['class'] = train_y
4 mean_ = np.array(df.groupby('class').mean())
5 for i in range(len(mean_)):
6     image = mean_[i]
7     plt.imshow(image.reshape((28, 28)), cmap = 'Greys')
8     plt.show()
```

```
1 ###compare with library build-in predict function
2 dataset = [(train_x_thre, test_x_thre, 'untouched'), (train_x_crop_thre, test_x_crop_thre, 'stretched')]
3 from collections import defaultdict
4 result = defaultdict(list)
5 for (train, test, name) in dataset:
6     for model in [GaussianNB(), BernoulliNB()]:
7         model.fit(train, train_y)
8         pred_y = model.predict(test)
9         train_score = accuracy_score(map(int, train_y), model.predict(train))
10        test_score = accuracy_score(map(int, test_y), pred_y)
11        result[name].append((train_score, test_score))
12 print(result)
13 ### result {'untouched': [(Gaussian_train_accu, Gaussian_test_accu),
14 ###                        (Bernoulli_train_accu, Bernoulli_test_accu)],
15 ###        'stretched': [(Gaussian_train_accu, Gaussian_test_accu),
16 ###                      (Bernoulli_train_accu, Bernoulli_test_accu)]}
```

```
defaultdict(<type 'list'>, {'untouched': [(0.54855, 0.5399), (0.8376833333333333, 0.8487)], 'stretched': [(0.79525, 0.8075), (0.8250333333333333, 0.8381)]})
```

Part 2B: MNIST using Decision Forest

```
1 from sklearn.ensemble import RandomForestClassifier
```

```
1 import itertools
2 trees = [10, 30]
3 depth = [4, 16]
4 combinations = list(itertools.product(trees, depth))
5 combinations
```

```
[(10, 4), (10, 16), (30, 4), (30, 16)]
```

```
1 resultB = defaultdict(list)
2 for (train, test, name) in dataset:
3     for [tree, depth] in combinations:
4         model = RandomForestClassifier(n_estimators = tree, max_depth = depth)
5         model.fit(train, train_y)
6         pred_y = model.predict(test)
7         train_score = accuracy_score(map(int, train_y), model.predict(train))
8         test_score = accuracy_score(map(int, test_y), pred_y)
9         resultB[name].append((train_score, test_score))
10 print(resultB)
11 ###resultB = {'untouched': [(train_accu 10trees + 4depth, test_accu 10trees + 4depth),
12 ###                        (train_accu 10trees + 16depth, test_accu 10trees + 16depth),
13 ###                        (train_accu 30trees + 4depth, test_accu 30trees + 4depth),
14 ###                        (train_accu 30trees + 16depth, test_accu 30trees + 4depth)],
15 ###        'stretched': [(train_accu 10trees + 4depth, test_accu 10trees + 4depth),
16 ###                      (train_accu 10trees + 16depth, test_accu 10trees + 16depth),
17 ###                      (train_accu 30trees + 4depth, test_accu 30trees + 4depth),
18 ###                      (train_accu 30trees + 16depth, test_accu 30trees + 4depth)]}
```

```
defaultdict(<type 'list'>, {'untouched': [(0.7433166666666666, 0.7615), (0.9946333333333334, 0.9489), (0.7983, 0.8061), (0.9976, 0.962)], 'stretched': [(0.7162166666666666, 0.7298), (0.9948333333333333, 0.9553), (0.742, 0.7576), (0.9708333333333333, 0.9672)]})
```